# Allinea DDT: Your Partner in Finding Debugged Paths on Mira

Ian Lumb <ilumb@allinea.com>
Senior Systems Engineer, Allinea Software Inc.

Mira Community Conference 2013

# [L2P] Summary



Leap to Petascale
Making the Move Towards Mira
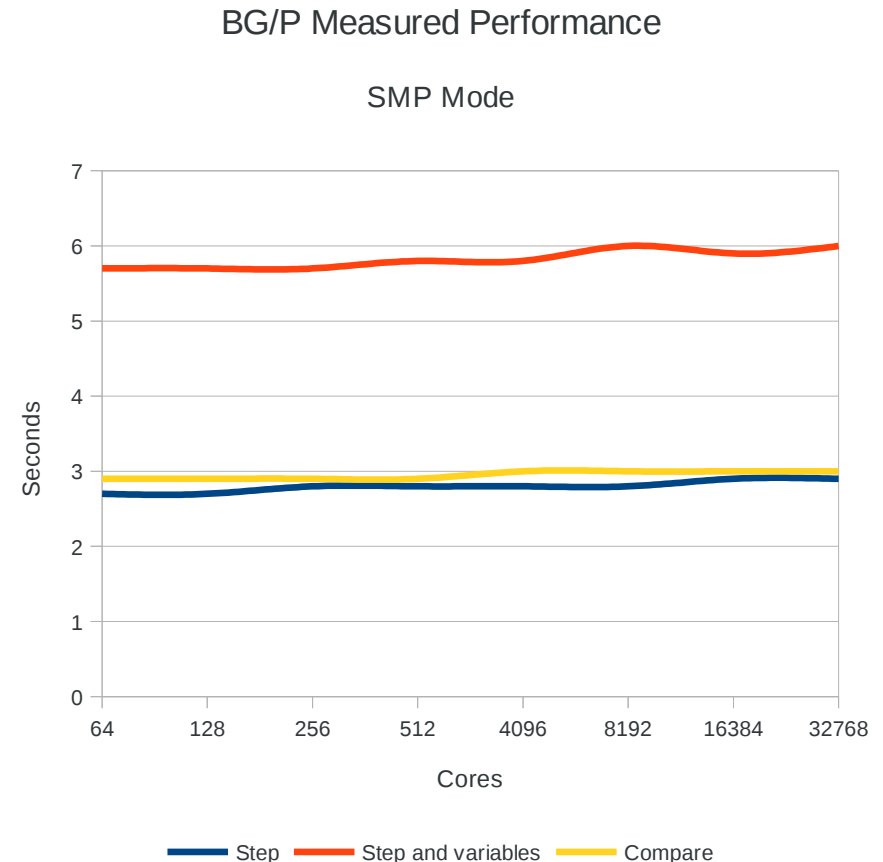May 22-25 -Argonne National Laboratory

- Petascaling for > 1 year
  - Petascaled infrastructure and UI
- Scaling for IBM Blue Gene /P
  - Acceptance testing at ALCF
- Scaling for IBM Blue Gene /Q
  - Addressing ALCF requirements
    - Early access for IBM Blue Gene /Q expected July 2012
- Architecture applicable elsewhere
  - Multicore/GPU??? architectures
- Exascaling ...

http://www.alcf.anl.gov/sites/www.alcf.anl.gov/files/L2PAllinea_0.pdf

allinea
www.allinea.com

# A Path to Petascale on IBM BG /P

- **Phase 1 [2010]**

  - Cut memory usage per compute process at I/O node

  - Debuggers share common internal tables

    - Memory mapping of symbol tables
    - Raises limit to ~128 processes
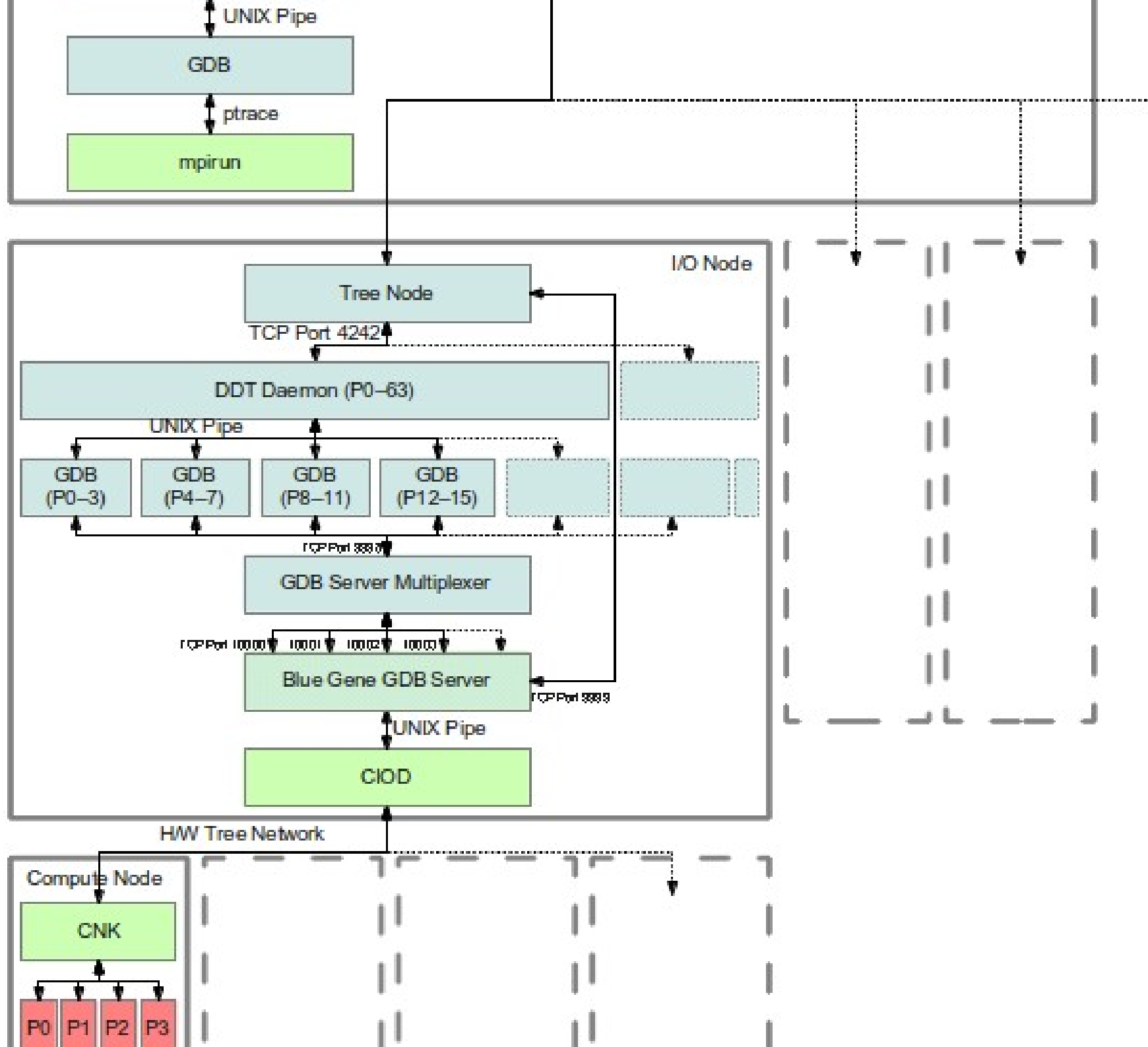
- **Delivered!**

# The memory mapped result

- Simplest to achieve – with benefits to multicore systems

  - Boosted max cores per I/O node to 256

- Reached 32K cores

  - 32,000 cores as quick as 64 cores

  - … flat – but not instantaneous

  - Most operations ~ 3 seconds

  - Close work with ANL – ran at scale on Intrepid

BG/P Measured Performance

SMP Mode



Step | Step and variables | Compare

www.allinea.com

# Petascale IBM Blue Gene /P Debugging

- Phase 2 [2011]
  - Reduce per-I/O-node daemon count
  - Reduces context thrashing: faster!
  - Each daemon handles multiple compute processes
    - Multiplexing commands and responses via CIOD
    - Multiplexing within the debugger
    - Cuts memory usage and improves speed
  - Limit 256-512 processes per I/O node
- Delivery: July 2012

**Multiplexed Architecture**

UNIX Pipe

GDB

ptrace

mpirun

I/O Node

Tree Node

TCP Port 4242

DDT Daemon (P0–63)

UNIX Pipe

| GDB (P0–3) | GDB (P4–7) | GDB (P8–11) | GDB (P12–15) |

TCP Port 8000

GDB Server Multiplexer

TCP Port 10000   10001   10002   10003

Blue Gene GDB Server

TCP Port 8889

UNIX Pipe

CIOD

HW Tree Network

Compute Node

CNK

| P0 | P1 | P2 | P3 |

# BG /P Case Study: Background

- Outstanding problems in heliospheric physics

  - Origin of the solar wind

  - Heating of the solar corona

- Large-scale numerical simulations

  - ***Simulation crashes at 16,386 MPI processes***

# Why debug at scale?

- **Increasing job sizes leads to unanticipated errors**

  - **Regular bugs**

    - Logic issues and control flow

    - Data issues from larger data sets – eg. garbage in..., overflow

  - **Increasing probability of independent random error**

    - Memory errors/exhaustion – "random" bugs!

    - System problems – MPI and operating system

  - **Coded boundaries**

    - Algorithmic (performance) or hard-wired limits ("magic numbers")

  - **Unknown unknowns**

- **Machine time is too expensive to ignore failures!**

allinea
www.allinea.com

- ## Reproduced the crash

  - ### Ran Allinea DDT in offline mode

    - #### Viewed HTML results via Web browser

      - Crash inside an MPI function call on about 128 of the 16384 cores
        - MPI implementation bug?
        - Memory bug?

- ## Ran Allinea DDT in offline mode again

  - ### Memory debugging enabled

    - #### Crash inside a harmless looking loop

      - Issue with loop index

- ## Ran Allinea DDT in GUI mode

  - ### Early calculation of the X-Y-Z grid is incorrect

allinea

# TRAFFIC

- Debugging
  - Transforming a broken program into a working one
- **How?**

  - **T**rack the problem

  - **R**eproduce

  - **A**utomate - (and simplify) the test case

  - **F**ind origins – where could the "infection" be from?

  - *Focus – examine the origins*

  - *Isolate – narrow down the origins*

  - *Correct – fix and verify the testcase is successful*

Zeller A., "Why Programs Fail", 2nd Edition, 2009

# Allinea DDT 3.2.1 – October 2012

# Allinea DDT and Mira

"This tool has already proven its value in the migration of our early science applications onto Mira," said Kalyan Kumaran, who manages ALCF's applications performance engineering team. "These projects cover the range of scientific fields, numerical methods, programming models and computational approaches expected to run on Mira, so **accurate debugging is critical**."

# Allinea DDT
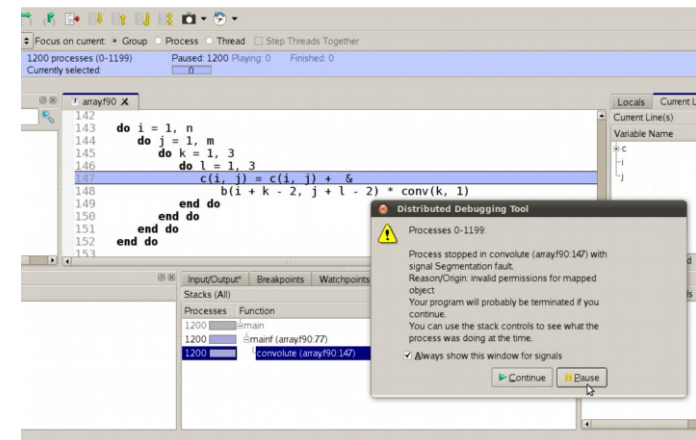
- Licensing
  - 32K-process permanent license
  - Full machine development license available (contact support)
- Startup overview
  - Compile –g –O0
    - OMP code compile -qsmp=omp:noauto:noopt
  - Softenv key "+ddt"
  - Need X11 server and ssh –X forwarding
  - [BG/P only] Start interactive job with *isub*
  - [BG/P or BG/Q] Run ddt and submit job through GUI
- More details:
  - [BG/P] http://www.alcf.anl.gov/resource-guides/allinea-ddt

http://www.alcf.anl.gov/sites/www.alcf.anl.gov/files/2013MiraCon_debugging_0.pdf
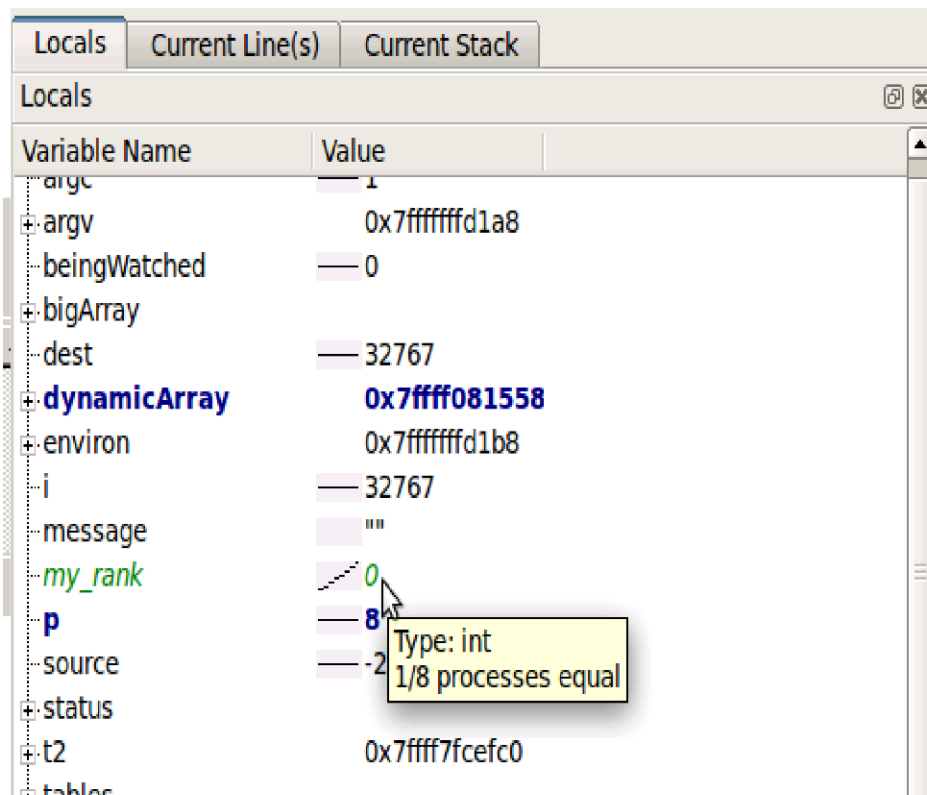
# Fixing the everyday crash

- The typical application crash or early exit:
    - Run your program in the debugger

        ddt {application} {parameters}
    - Application crashes or starts to exit

- **Where** did it happen?
    - Allinea DDT merges stacks from processes and threads into a tree
    - Leaps to source automatically

- **Why** did it happen?
    - Some faults evident instantly
    - For others look deeper – at variables

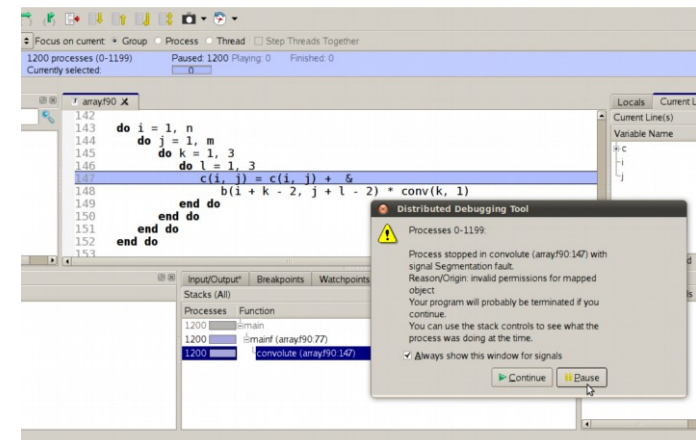# Simplifying the data deluge



- Allinea DDT compares data automatically
  - Too many variables to trawl manually!
- Smart highlighting
  - Subtle hints for differences and changes
  - With sparklines!

- More detailed analysis
  - Full cross process comparison
  - Historical values via tracepoints

allinea
www.allinea.com

# Allinea DDT: Proved to the extreme

- Scalability by design
  - User interface that scales
  - High performance tree architecture
- Proven performance at Petascale
  - Measured in milliseconds
  - **Routine use** at 100,000+ cores
- 300,000+ cores
  - Easy to use
  - Scalable GUI

# Allinea DDT: More than debugger

- **Integrated automated detection of bugs**
  - Static analysis
  - Memory leaks and errors
- **Open plugin architecture**
  - MPI checking tools
- **Offline mode - debug in batch mode**

# Allinea DDT - Debugging++

- Productively **debug** your parallel code

- Completely **understand** your parallel code

  - Interact with data, algorithms, codes, programs and applications in real time

- **Develop** parallel your code from scratch

- **Port** parallel algorithms, codes, programs and applications to X

- **Scale** your algorithms, codes, programs and applications

# The Allinea Environment: Benefits

- **At last:** a modern **integrated** environment for the HPC developer

- Supporting the lifecycle of application development and improvement
  - Productively debug code
  - Enhance application performance

- Designed for productivity
  - Consistent integrated easy to use tools
  - Enables effective HPC development

- Improve system usage
  - Fewer failed jobs
  - Higher application performance



PERFORM

the
**allinea**
environment

FIX

**allinea**
www.allinea.com

# What's really new?



STAY TUNED

allinea
www.allinea.com